# Index

# Topics, goals and urls

- Is there any separability in how much you have been liked among the participants?
- Which attributes of a partner do participants attach most importance to?

https://www.kaggle.com/datasets/ulrikthygepedersen/speed-dating

# Data mining process

# Data cleaning

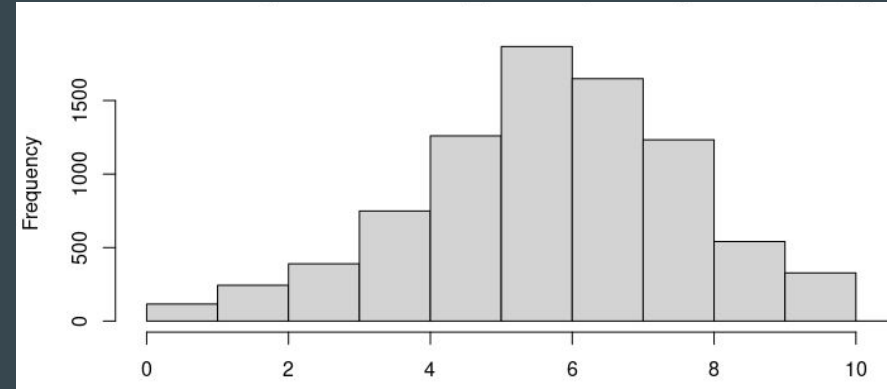| Reducing the size of our dataset | Changing the format of the values | Erase variables with unacceptable number of null values |
|:---:|:---:|:---:|
| | d'0' → 0 | |
| | d'1' → 1 | |
| 8378 x 123 | | expected_num_interested_in_me |
| ↓ | [0...1] → 1 | |
| | [1...2] → 2 | **79%** NULL |
| | [2...3] → 3 | |
| 8378 x 25 | . | |
| | . | |
| | . | |
| | [9...10] → 10 | |

# Preprocessing - Imputation of null values

## Numerical variables: Mean Imputation

Histogram without imputation

Histogram with mean imputation



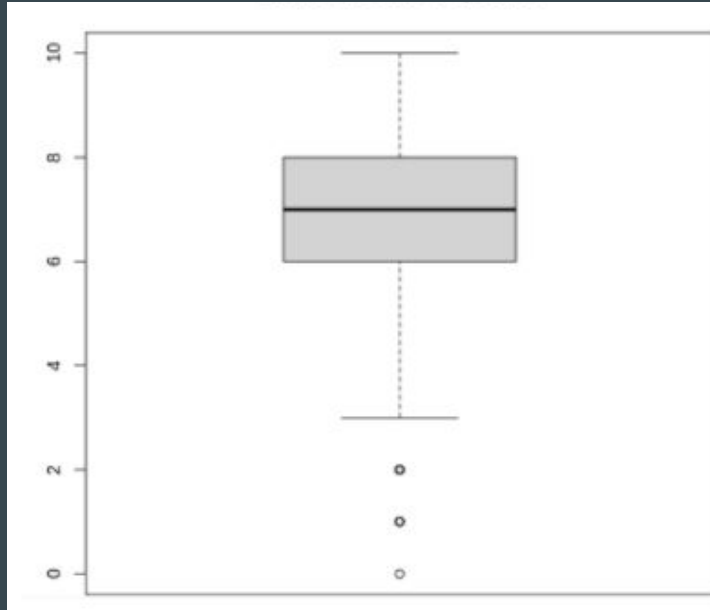## Categorical variables: MCA Imputation

# Preprocessing - Errors in values

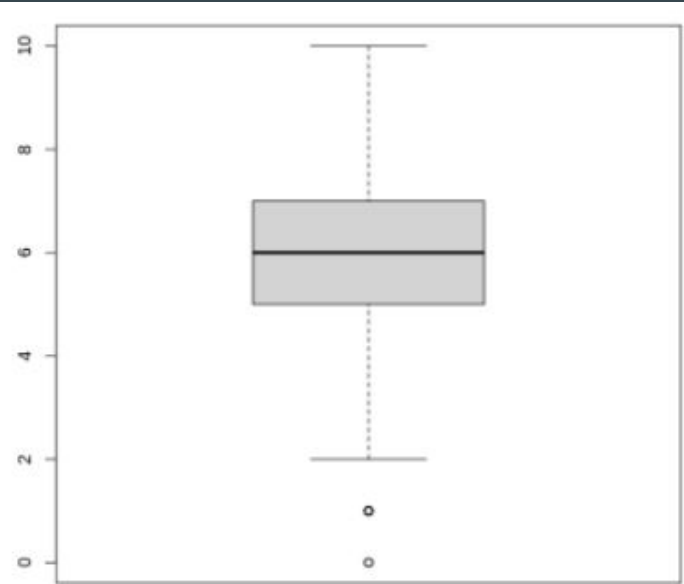- Are the variables in the range that they should be according to the metadatafile?

```
attractive_o
Min.    : 0.000
1st Qu.: 5.000
Median : 6.000
Mean   : 6.186
3rd Qu.: 8.000
Max.   :10.500
```

# Preprocessing - Outliers treatment

BoxPlot "sincere_partner"

BoxPlot "like"

# Descriptive analysis

## Gender at birth



**Pie of gender**

Binary

## Age difference



**Histogram of d_age**

Numerical

# Descriptive analysis

Interest correlation



Importance of attractiveness



Numerical

Categorical

# PCA

# PCA

Projection of numerical variables



**Projection of numeric variables**

# PCA

Subspace
(1,3)

Projection of numerical
and categorical variables



**Projection of categories**
ref_o_sincere, pref_o_intelligence, attractive_important, sincere_importan

# Clustering  (only numerical variables)

Dendogram with distance

Dendogram with distMatrix



196 individuals of 8378 in c2

4112 individuals of 8378 in c2

# Clustering

# Profiling



Boxplot of d_age vs Cluster

Means of d_age by Cluster

Boxplot of attractive_o vs Cluster



Means of attractive_o by Cluster

**Boxplot of sinsere_o vs Cluster**

**Means of sinsere_o by Cluster**

Boxplot of intelligence_o vs Cluster

Means of intelligence_o by Cluster

Prop. of pos & neg by attractive_important

Boxplot of like vs Cluster

Means of like by Cluster

# Conclusions



- Separability in how much their partners like them

- Intelligence > Sincere > Attractive

# Original Scheduling

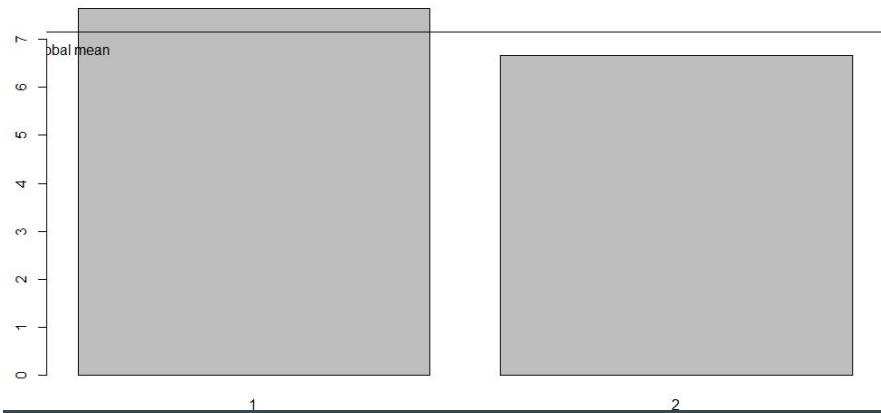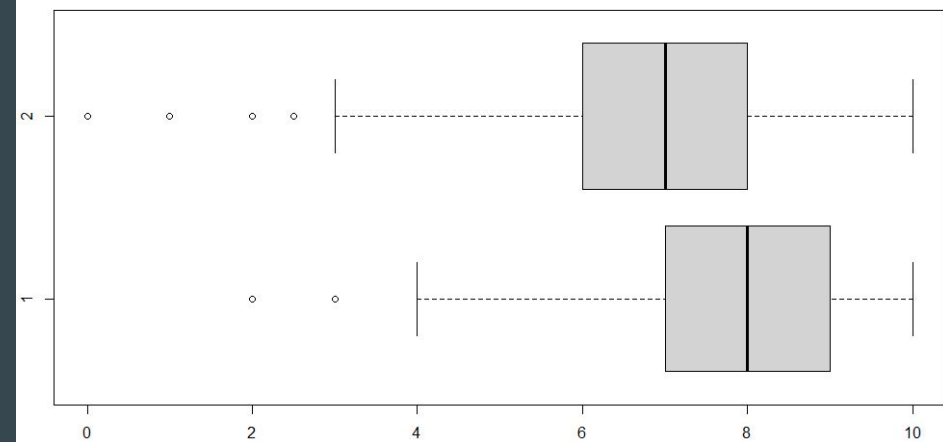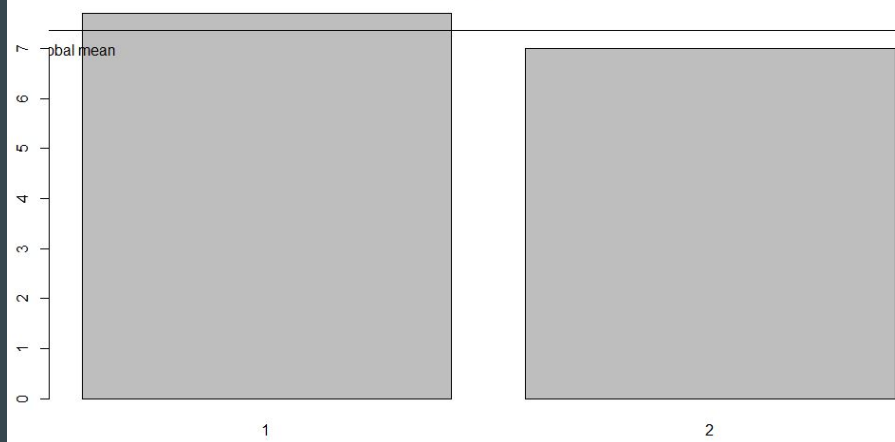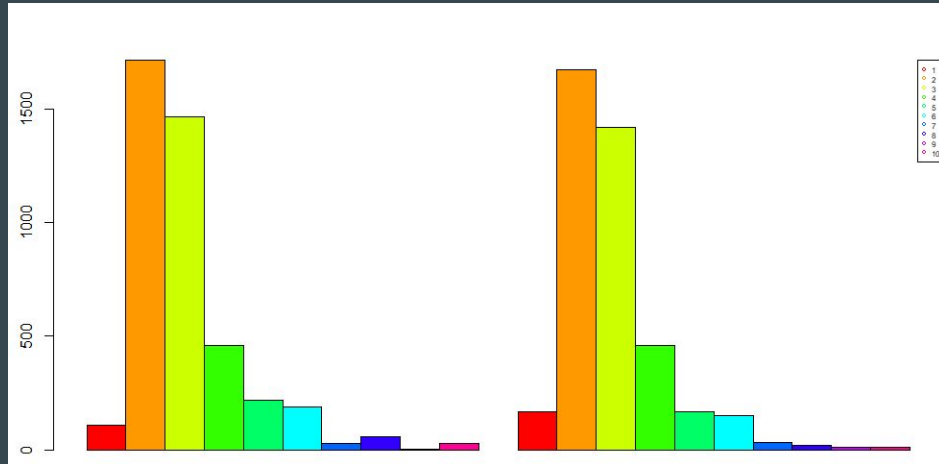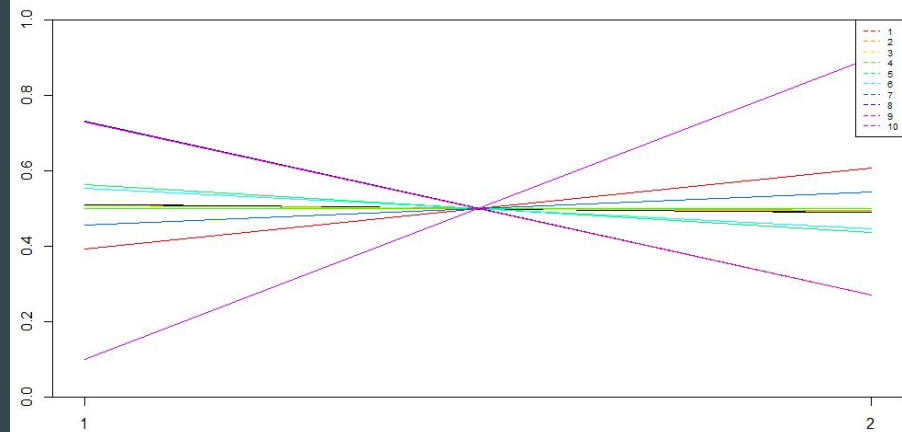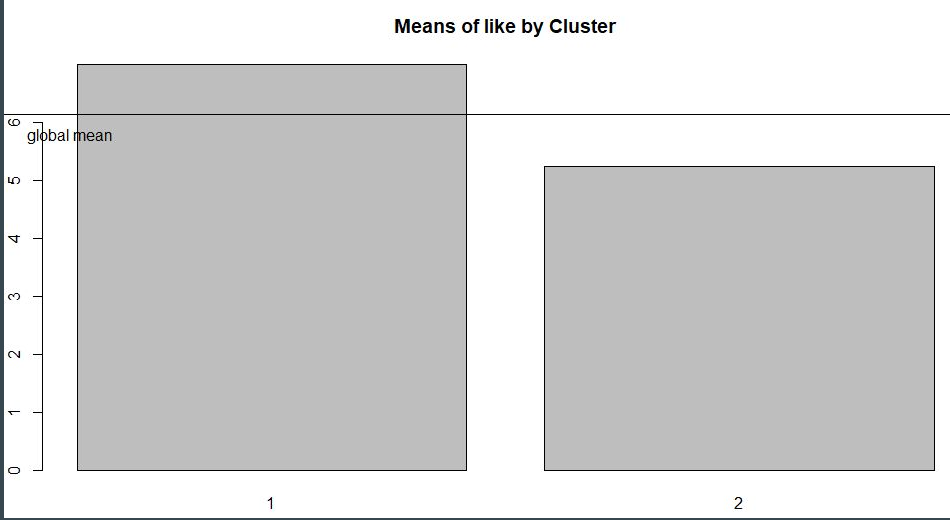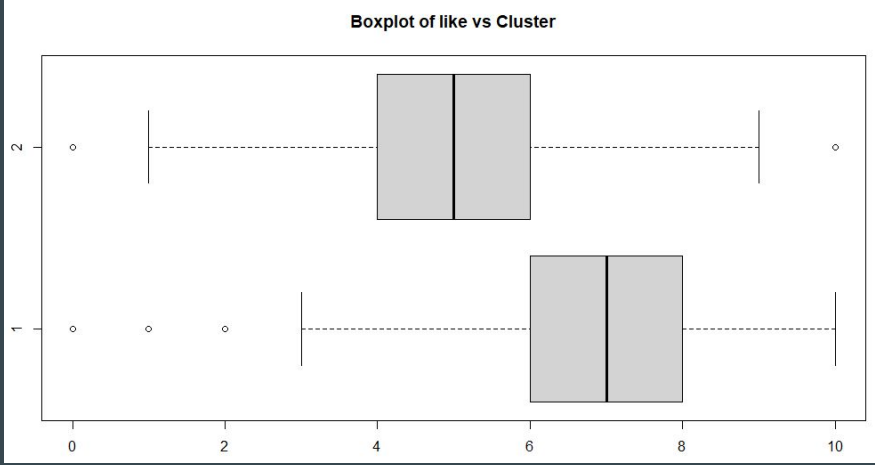| Tasks | February 24 | 25 | 26 | 27 | 28 | March 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Assignment Grid | ■ | ■ | ■ | ■ | ■ | ■ | ■ | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Gantt Chart | ■ | ■ | ■ | ■ | ■ | ■ | ■ | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Risk Plan | ■ | ■ | ■ | ■ | ■ | ■ | ■ | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Metadata file | | | | | | | | ■ | ■ | ■ | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Description of raw variables | | | | | | | ■ | ■ | ■ | ■ | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Preprocessing steps | | | | | | | | | | | ■ | ■ | ■ | ■ | ■ | ■ | | | | | | | | | | | | | | | | | | | | |
| List and justify preprocessing decisions | | | | | | | | | | | | | | ■ | ■ | ■ | ■ | ■ | | | | | | | | | | | | | | | | | | |
| Descriptive statistics of modified variables | | | | | | | | | | | | | | | | | | | ■ | ■ | ■ | | | | | | | | | | | | | | | |
| PCA | | | | | | | | | | | | | | | | | | | | | | ■ | ■ | ■ | | | | | | | | | | | | |
| Herarchial Clustering | | | | | | | | | | | | | | | | | | | | | | | | | ■ | ■ | ■ | ■ | | | | | | | | |
| Profiling of clusters | | | | | | | | | | | | | | | | | | | | | | | | | | | | | ■ | ■ | | | | | | |
| Conclusions | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | ■ | ■ | | | | |

# Final Scheduling

| Tasks | February 24 | 25 | 26 | 27 | 28 | March 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Assignment Grid | ■ | ■ | ■ | ■ | ■ | ■ | ■ | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Gantt Chart | ■ | ■ | ■ | ■ | ■ | ■ | ■ | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Risk Plan | ■ | ■ | ■ | ■ | ■ | ■ | ■ | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Metadata file | | | | | | | | ■ | ■ | ■ | ■ | | | | | | | | | | | | | | | | | | | | | | | | | |
| Description of raw variables | | | | | | | | ■ | ■ | ■ | ■ | | | | | | | | | | | | | | | | | | | | | | | | | |
| Preprocessing steps | | | | | | | | | | | ■ | ■ | ■ | ■ | ■ | ■ | | | | | | | | | | | | | | | | | | | | |
| List and justify preprocessing decisions | | | | | | | | | | | | | | ■ | ■ | ■ | ■ | | | | | | | | | | | | | | | | | | | |
| Descriptive statistics of modified variables | | | | | | | | | | | | | | | | | | | ■ | ■ | ■ | ■ | | | | | | | | | | | | | | |
| PCA | | | | | | | | | | | | | | | | | | | | | | ■ | ■ | | | | | | | | | | | | | |
| Herarchial Clustering | | | | | | | | | | | | | | | | | | | | | | | ■ | ■ | | | | | | | | | | | | |
| Profiling of clusters | | | | | | | | | | | | | | | | | | | | | | | | | ■ | ■ | ■ | | | | | | | | | |
| Conclusions | | | | | | | | | | | | | | | | | | | | | | | | | | | | ■ | ■ | ■ | ■ | | | | | |

# Questions ?